

UNIT - 02

What is ETL - ETL vs ELT - Types of Data Warehouse -
 DW design & Modeling - Delivery process - OLAP [Online Analytical processing] - Characteristics of OLAP -
 OLTP [Online Transaction processing] vs OLAP - OLAP Operations - Types of OLAP - ROLAP vs MOLAP vs HOLAP

TOPIC NAME :- What is ETL :-

- Extract Transform/Load
- Process of combining multiple source called DW
- Organize raw data & Machine Learning (ML)

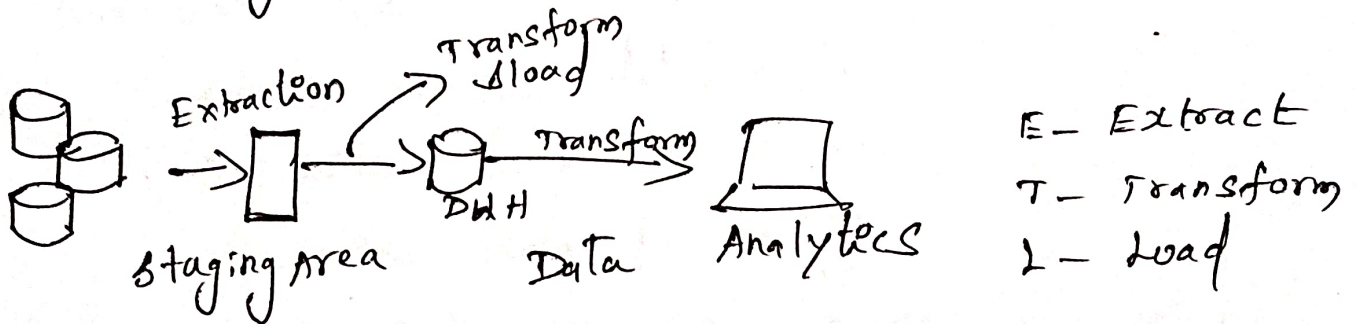


Fig :- ETL

Why ETL is important :-

- Both structure & unstructure data
- CRM - customer Relationship management
- IIOT - Internet of things
- Prepare structure for analytic purpose
- Study customer behaviour

ETL Evaluation:-

Convert transactions to relational data formats for analysis

Traditional ETL:-

- support read & write request
- analyze popular items
- Identify relationship b/w table / pattern / trends

Modern ETL:-

create data sinks & receive data from multiple sources

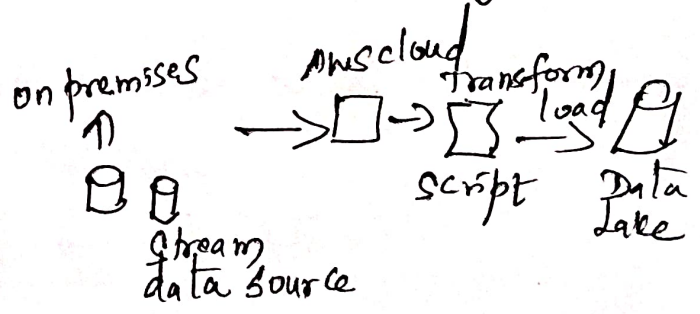
Data Ware houses:-

- store multiple database
- organize table & column
- optimize data processing

Data lakes:- Data like SQL queries / ML / Big data etc

How does ETL work:-

- source db
- Transform data
- load data



EXTRACTION:-

dig:- IETL process

- Extract raw data from multiple resources
- happen
 - update notification - Extract changes
 - Incremental Extraction - notify data record
 - Full extraction - high data transfers

TRANSFORMATION:-

- consolidate raw data & involve
- Basic data transformation - simplify data
- Data cleansing - remove error
- Data deduplication - remove duplicate record
- data format session - convert into pounds
- Advanced data transformation - easy analysis
- Deriving - calculate new value from existing value
- Joining - join link & data sources
- splitting - divide target system
- Summarization - reduce data volume
- Encryption - provide data stream

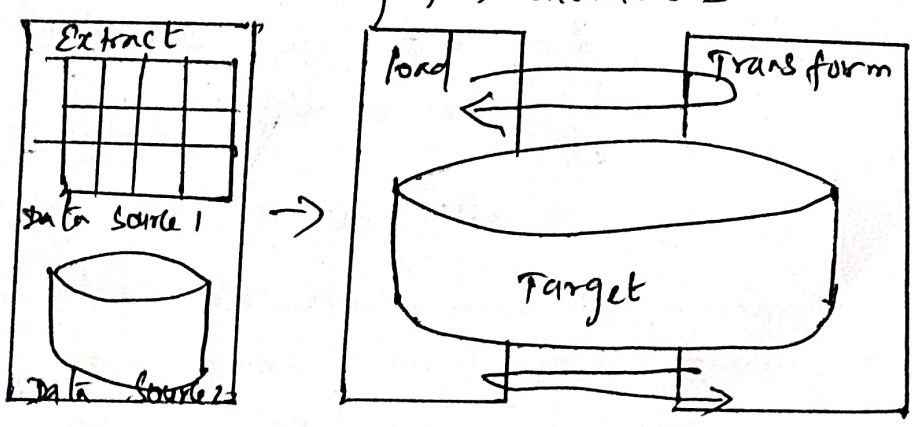
LOADING :- Well defined

- full load - load data
- Incremental load
 - Streaming Incremental Load
 - Batch Incremental load

What is ETL:-

Extract, Load & transform

ETL process - extract data & needs to be loaded repository as well as perform calculations



ETL process

→ Benefits of ELT:-

- Extract, Transform, Load
- Separate transformation step
- Save time & resources
- Easy to change as per requirement

→ ETL VS ELT:-

ELT load raw data to target
 ETL - transforms secondary processing server

→ Data compatibility:-

represent table rows & columns
 handle all type of data like document/images

→ Speed:-

difficult to scale / data size increases

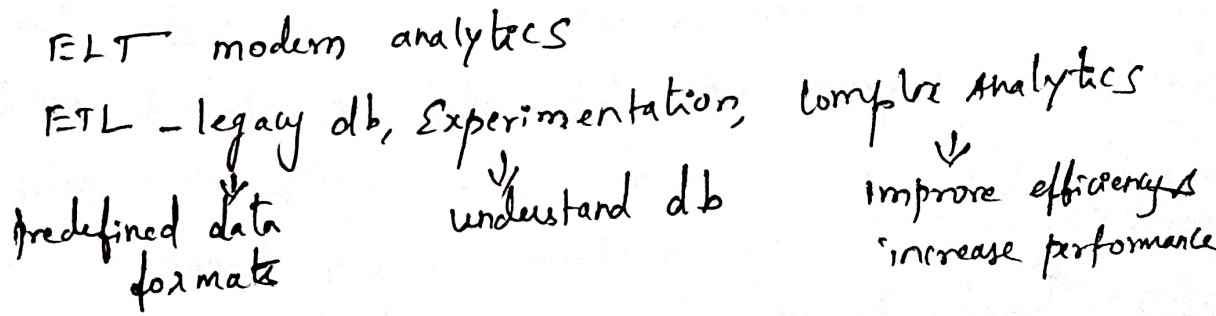
→ Cost:-

More, simpler data stack & lower setup costs

→ Security:-

protect data

When to use ETL VS ELT:-



IOT Applications:- use sensor data streams
 filtered value reduced

Differences:-

A) ETL VS ELT process order

ELT - data extracted from its source
 more flexible transformation & easily modified
 ETL efficient approach & process target system

- B) ETL VS ELT - MAINTENANCE
 run efficiently for monitoring / scheduling / Backup / data quality
 ensure target system for data transformation
- C) ETL VS ELT : COST
 data egress cost / less data transforms
- D) ETL VS ELT : Security
 perform separate queries
 access only specific sources
 less secure
- E) ETL VS ELT : Hardware requirements
 powerful & robust
 handle large volume of data
 manage data as specialized ETL tool used
- F) ETL VS ELT : Support for Data Lake
 hold large amount of structured data
 cost effective, more complex setup
 better suited for lake data environment
 allow more flexible & scalable
- G) ETL VS ELT - Support for Data Warehouse
 deal less data cloud lakes as properly structured one
- H) ETL VS ELT - performance
 ensure large data volumes as data transformation is parallel
 prefer less data size
- I) ETL VS ELT : Data movement
 intensive process, reduce amount of data

I) ETL VS ETL: FLEXIBILITY & SCALABILITY

⑤^{II}

Traditional data processing method

Easy to scale, depend on specific need of requirement

J) ETL VS ELT: SUPPORT FOR UNSTRUCTURED DATA

ETL/ELT tools used in business process

Incorporate user requirements

K) ETL VS ELT: Data Latency

Extract source data, result in higher data latency

for destination lower data latency occur.

Efficient but affect overall data latency of both process

L) ETL VS ELT: DATA LOSS:-

Cause data loss

ETL move all data into target systems

Storage requirements may be higher

M) ETL VS ELT: MATURITY:- Extract data for ETL process

CATEGORY	ETL	ELT
→ Stands for	extract/transform/load	Extract, load & transform
→ process	load predefined format	load target data
→ T & L location	Occurs in secondary process	take place target data
→ Data compatibility	Structured data	Structure/unstructured/ semi structured data
→ Speed	slower	faster
→ Cost	depend on ETL tool	depend on ELT used
→ Security	custom applications protection	manage data protection

TOPIC NAME:- TYPES OF DATA WAREHOUSES

- A) HOST BASED DW
- B) LAN BASED WORK GROUP
- C) SINGLE STAGE DATA WAREHOUSES
- D) MULTI STAGE DATA WAREHOUSES
- E) STATIONARY DATA WAREHOUSES
- F) DISTRIBUTED DATA WAREHOUSES
- G) VIRTUAL DATA WAREHOUSES

A) HOST BASED DW:-

Support robust / reliable
 handle centrally / Work group environment
 Allow automated extraction & cleaning of data
 Enable direct access
 Transaction processing oriented applications

Host-based (MVS) DW:-

- reside large volume of data
- high volume data storage

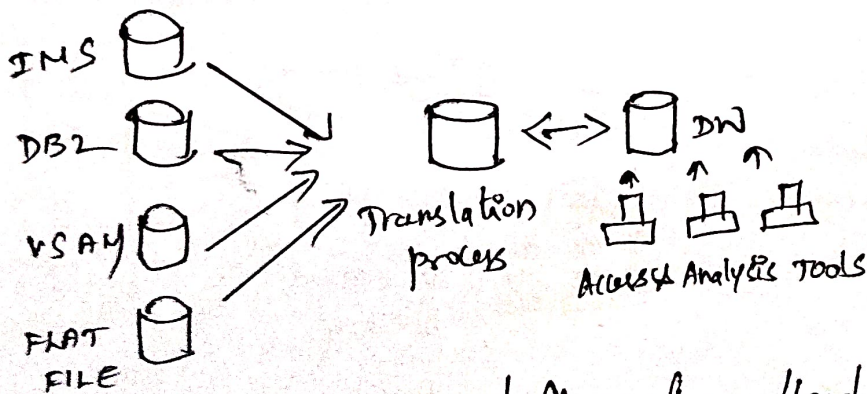


fig: Host based MVS/DW

- follow unload / transform / load phase
- design, build & maintain / end user access access technique
- support query / report / maintenance for DB2 tool

Host-Based UNIX DATA WAREHOUSES

Extract information UNIX Based dW
Support interworking feature

B) LAN BASED WORK GROUP DW:-

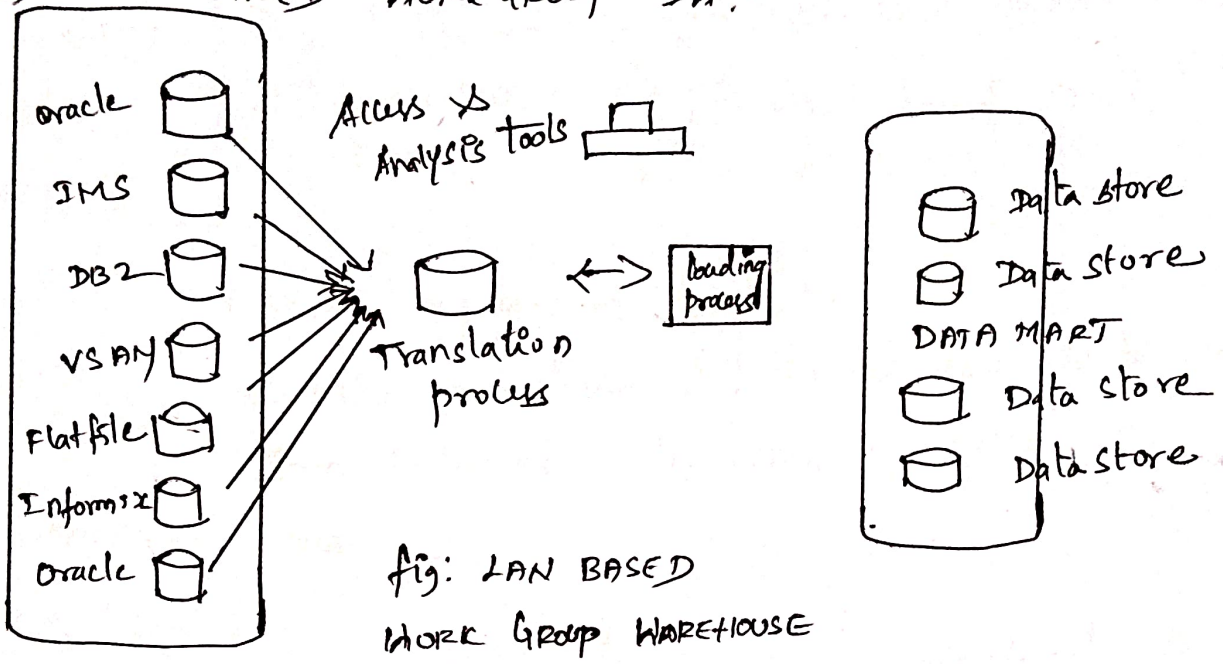


fig: LAN BASED WORK GROUP WAREHOUSE

Extract information multiple LAN
Data delivery - used in workgroup level
Maintain LAN environment / often called as data mart.

C) HOST BASED SINGLE STAGE (LAN) DATA WAREHOUSES:-

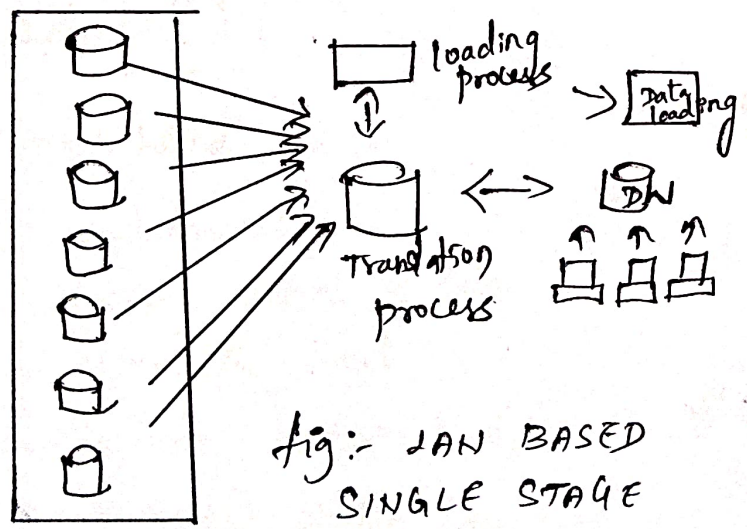
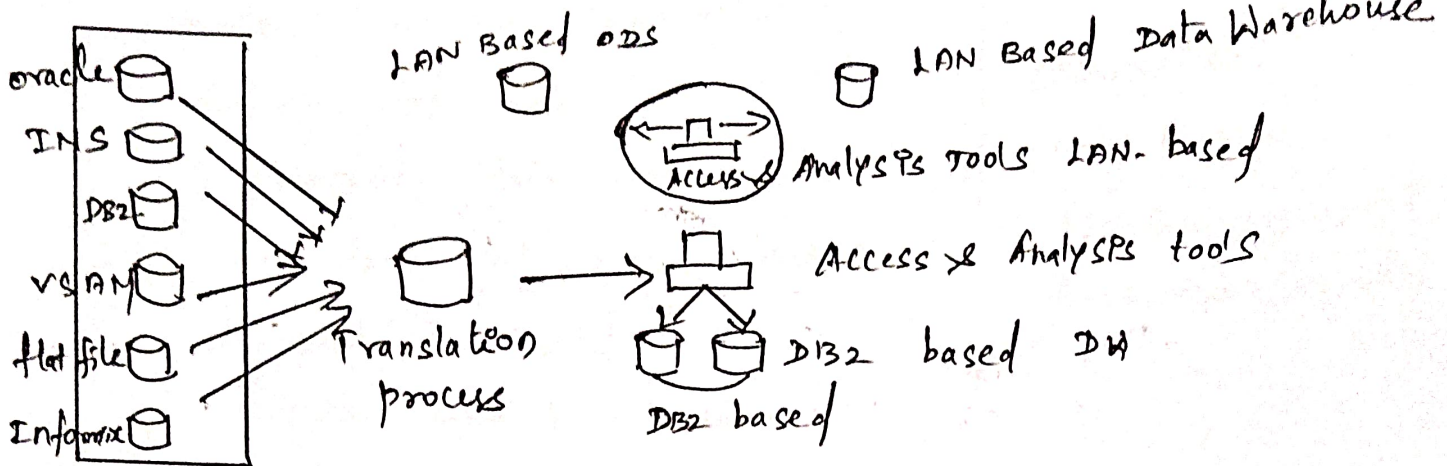


fig:- LAN BASED SINGLE STAGE WAREHOUSE

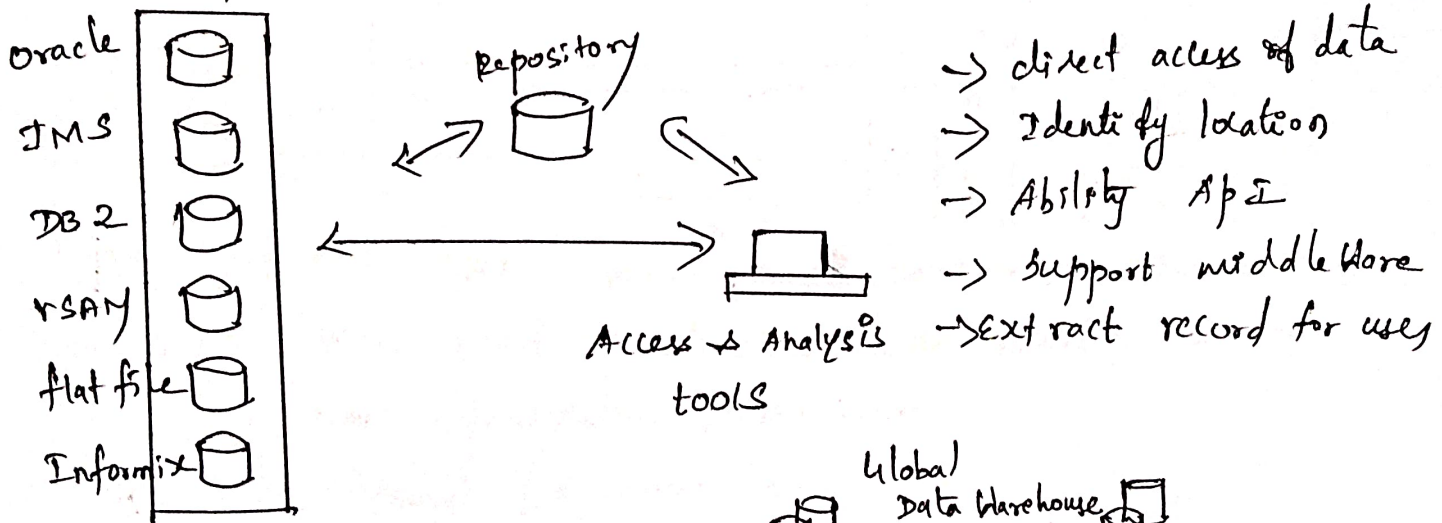
- handle Workgroup
- LAN based solutions
- dependent one
- design / plan / implementation
- support DB2 family
- support business data
- feasible

D) MULTI-STAGE DATA WAREHOUSES :-



support aggregation / data mart
 sustainable end client / historical data to be analyzed
 store historical calculation of files

E) STATIONARY DW :-



F) Distributed DW :-

- Autonomous
- unique content of data
- record is local
- Include historical data
- Integrate within local site

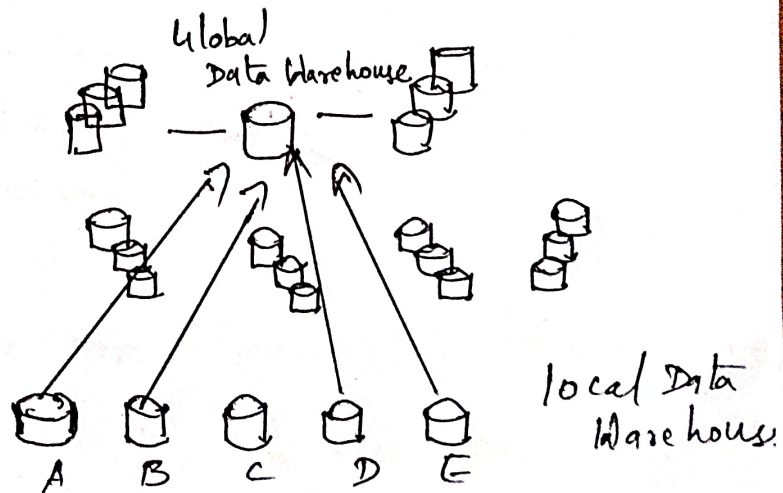


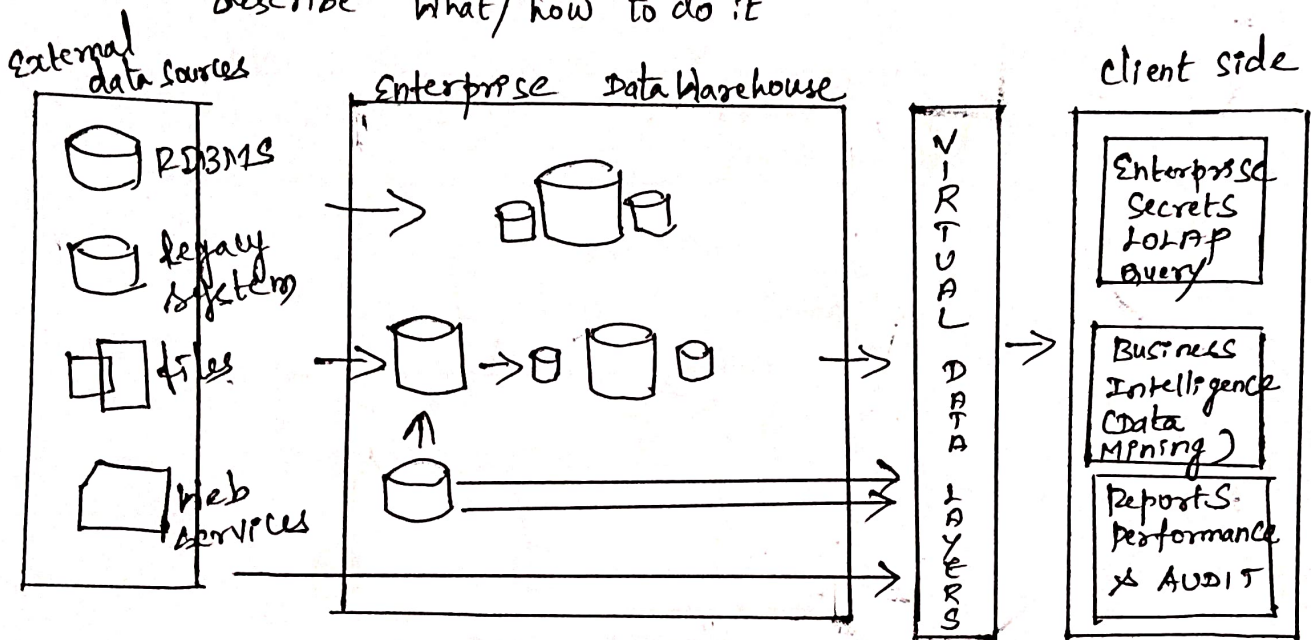
Fig: Distributed DW

G) VIRTUAL DW :-

10 "

process management facilities
 defers end users
 Implemented data access network
 redundant information loaded
 difficult to build
 describe what/how to do it

Disadvantage :-
 very complex
 No summary record
 performance degraded



TOPIC NAME: DATA WAREHOUSE DESIGN & MODELING

DESIGN :- OLAP used - online business analytical processing

Meet organization requirement

dynamic / continuous / ETL

approach 1) top down
 2) bottom up

1) Top Down design approach :-

Subject oriented / time variant / non volatile / data repository

stores atomic information / data must apply

single integrated data source used

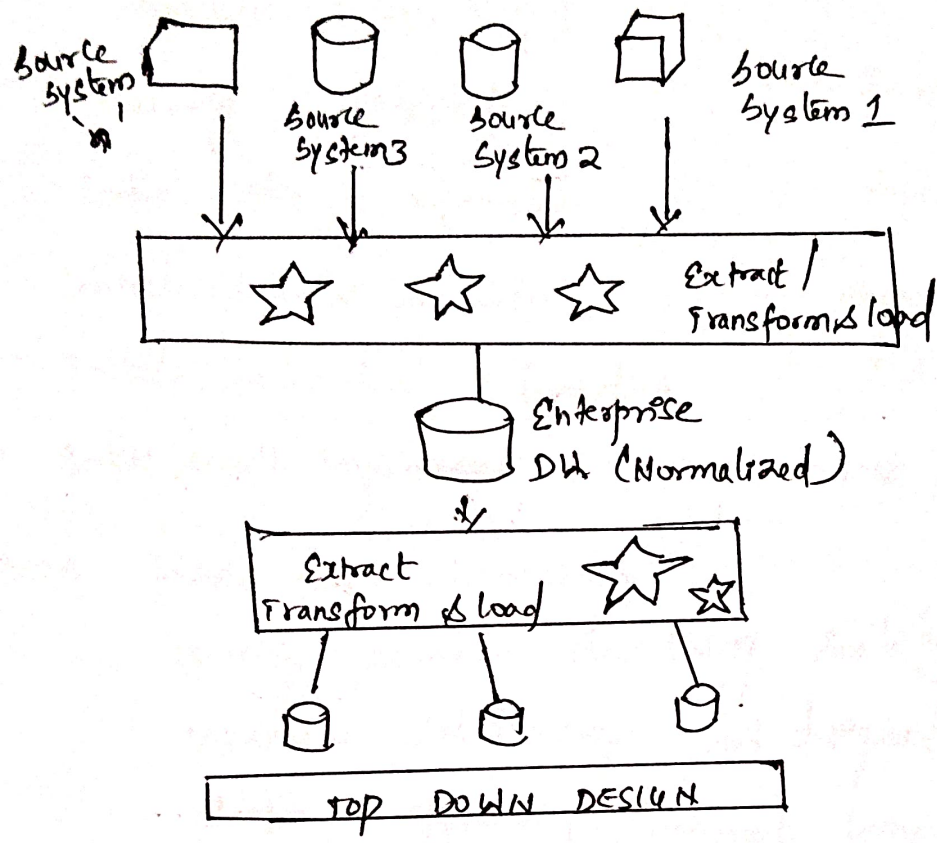
Advantage :-

Easy
Data mart loaded

disadvantage

Inflexible
Cost high

TOP DOWN DESIGN APPROACH

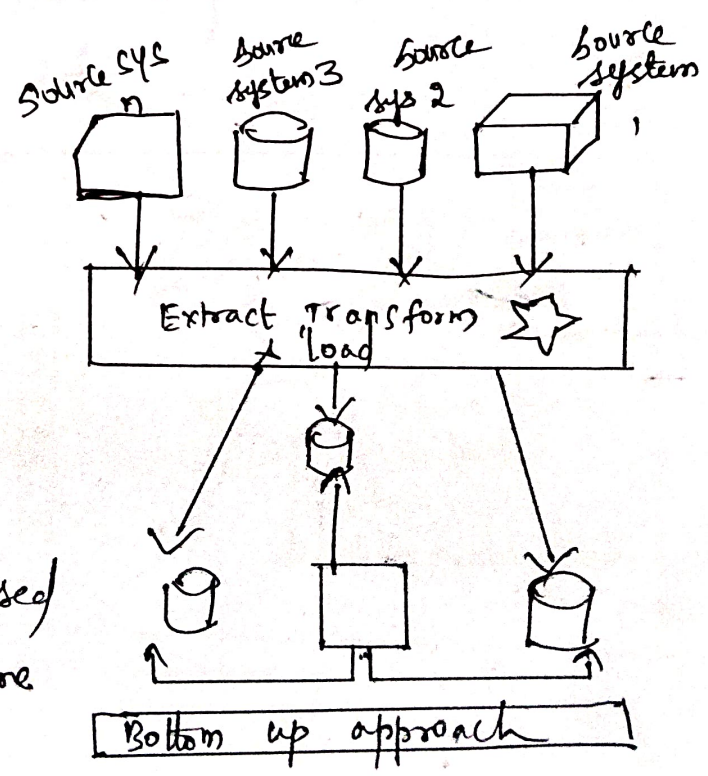


2) BOTTOM UP DESIGN APPROACH :-

query & analysis
business approach
virtual data warehouse
less time / failure less
learn & grow

Adv :- Generate quickly / Data mart used

disadv :- Bottom up approach - less failure

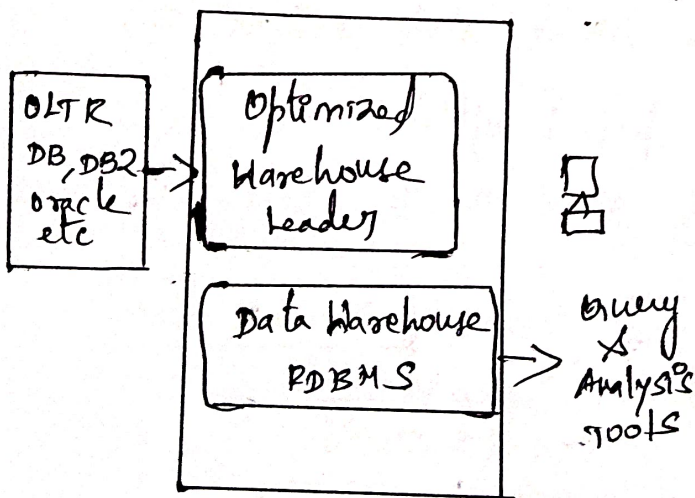


DIFFERENCE :-

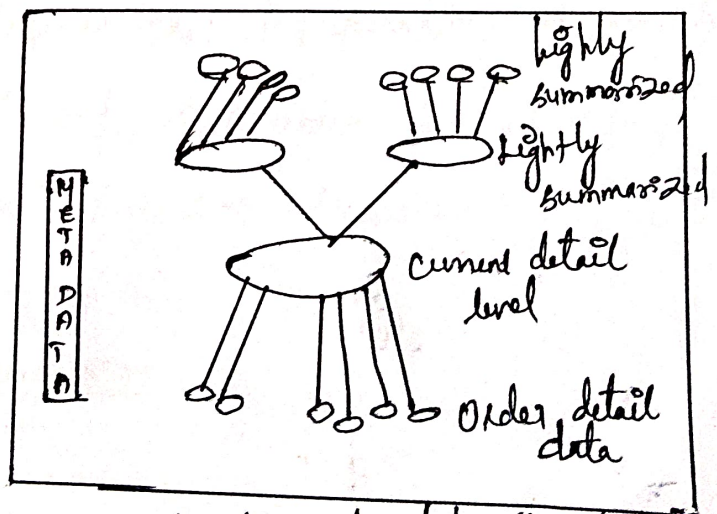
- | | |
|--|--|
| <p>Top Down design</p> <ul style="list-style-type: none"> - Breaks vast problem into smaller subproblems - Inherently architected - Central storage information - Centralized rules & control - Include redundant information - quick result implemented | <p>Bottom up design</p> <ul style="list-style-type: none"> - Integrate higher one - Inherently Incremental - departmental information stored - departmental rules & control - redundancy can be removed - less risk of failure |
|--|--|

DATA WAREHOUSE MODELING

develop schema, describing reality
 visualize relationship, well designed schema allow
 Improve efficiency / support complex queries
 Insert / delete / change data



DATA WAREHOUSE MODEL



structure of data inside the data Warehouse

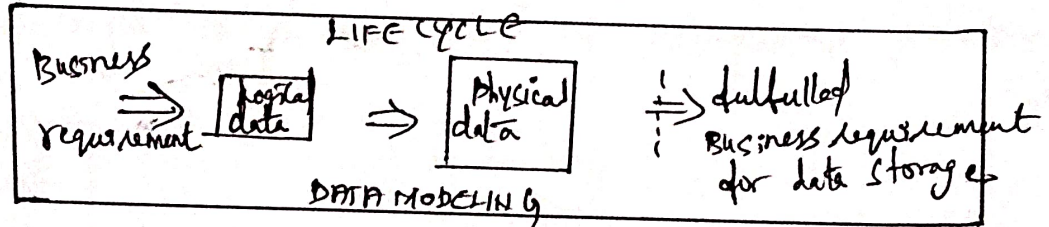
fast to access / expensive / difficult to manage
 order detail data - access current data

lightly summarized data - stored in disk storage / summarize attributes
 highly summarized data - outside DW

Metadata - locate data items / map record / summarization b/w current data

DATA MODELING LIFE CYCLE:-

straight forward process / fulfill goal for store, maintain, access



Conceptual Data Model:-

recognize relationship b/w different entities

Char: entities relationship

No attribute specified

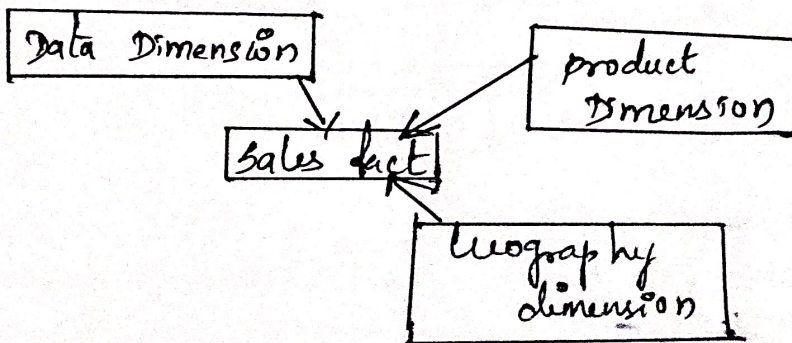
No primary key

Logical data model

→ logical data model defines the information as possible

→ No datatypes listed / referential Integrity specified

→ list all attribute for each entity / normalization.

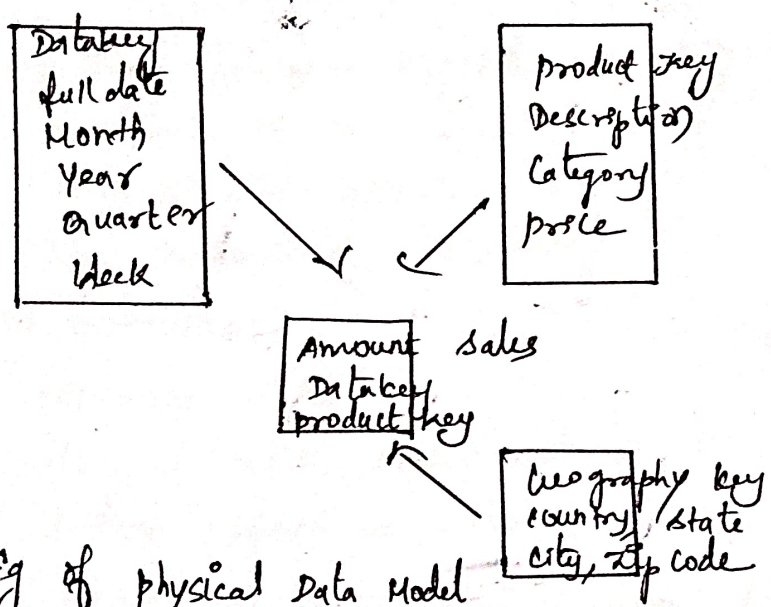


PHYSICAL DATA MODEL :-

- describe how the model is presented
- demonstrate table structures, column names, datatypes
- store information & include definition of new DS etc
- Enhancing query performance

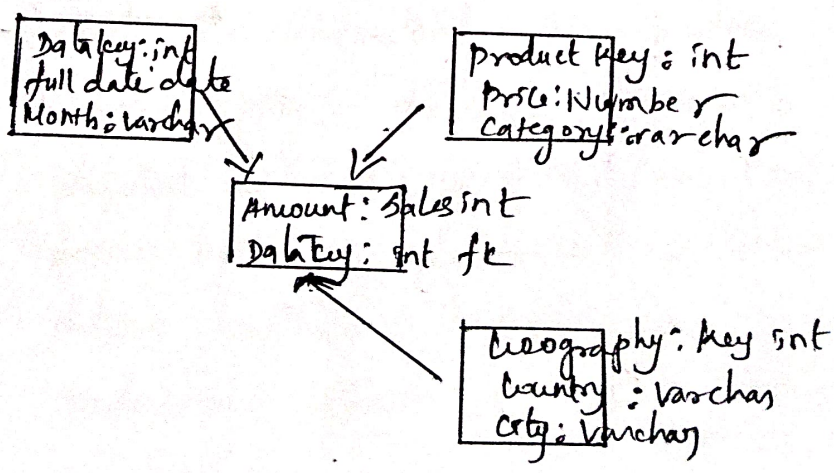
Characteristics of Physical Data Model :-

specification all tables & columns
foreign key are used to recognize relationship b/w tables



Eg of Physical Data Model

TYPES OF DATA WAREHOUSE MODELS :-



Enterprise Warehouse:-

- collect record, summarized detailed information
- support parallel architecture platforms
- used to develop & build

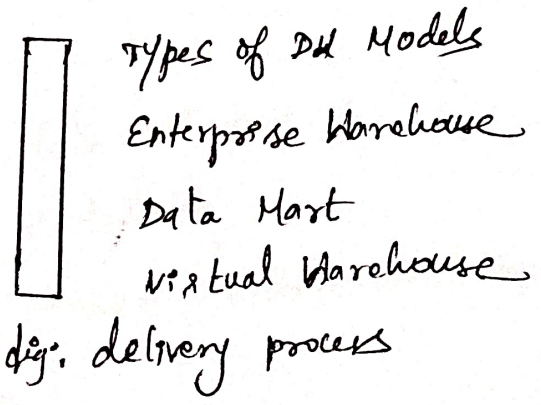
Data Mart:- divided into 2 parts

- Independent Data Mart
- Dependent Data Mart

VIRTUAL WAREHOUSES:- simple to build / support db servers

TOPIC NAME:- DELIVERY PROCESS:-

Evolve business expands
 designed to be flexible
 difficult to complete tasks
 Not understood completely



DELIVERY METHOD:-

- Minimize risks
- does not reduce overall delivery
- reduce project & delivery risk

IT strategy:- strategy retain funding for project

Business case: * tends to suffer delivery process

* needs to understand the business case for investment

Education & prototyping:- feasible & promote educational process

- defined technical objective / not critical
- focus business requirement & technical phase
- deliver business benefits
- understand short & medium term requirements

BUSINESS REQUIREMENTS:-

provide overall requirement
 source systems provide data
 fulfill short term requirement

Technical Blueprint:-

deliver overall architecture
 derive business benefits / data retention / backup & recovery

Building the version:-

add business benefit

History load:-

create to store / increased data volumes

allow user to analyze recent trends

Identify seasonal trends / complex / perform separate tasks

Adhoc query:- generate database query
 use access tools

automation:- fully automated, transforming data into analysis

* Backup, restore, archiving data

Extending scope:-

extend new set of business requirements

using existing information

Requirement Evolution:-

Not static / allow changes reflect within system

opposed to existing queries

operate pseudo application development process

rework overall system & continually update business needs

OLAP - ONLINE ANALYTICAL PROCESSING

OLAP is application architecture

Implement analytical application

using specialized MDDBS technology

Need for OLAP:-

array oriented / multidimensional nature

retrieve large number of records

represented relational db & accessed via SQL

Weakness series data & complex functions

OLAP is continuous, iterative & interactive process

OLAP Guidelines:-

does not eliminate methodology

Guidelines :- Multidimensional conceptual view - correspond business problem

Transparency - proficiency with front end environment tools

Accessibility - perform heterogeneous analysis

Consistent reporting performance - degradation in performance

client/server architecture - interoperability / flexibility / adaptivity

Generic dimensionality - structure & operational capabilities

Dynamic sparse matrix handling - achieve & maintain performance

Multisuser support - support work group of users

Unrestricted cross dimensional operations - perform rollup calculations

Intuitive data manipulation - drag & drop actions

flexible reporting - analyze analytical process

Unlimited dimensions & aggregation levels - OLAP not impose aggregate levels

TOPIC NAME: CHARACTERISTICS OF OLAP

A) MULTIDIMENSIONAL DATA ANALYSIS

3D graphics, advance computation & modeling & aggregation functions supports

B) ADVANCED DATABASE SUPPORT:-

Internal & external data sources / drill down & rollups support very large database

C) EASY TO USE END USER INTERFACES

GUI is simple / useful access

D) SUPPORT CLIENT / SERVER ARCHITECTURE

framework designed / developed / implemented divide OLAP system / distributed

BASIS	ROLAP	MOLAP	HOLAP
Storage	Relational db used	Multidimension db used	Multidimensional database used
Processing	ROLAP slow	MOLAP fast	HOLAP fast
Storage space	large	Medium	Small
Storage location	Relational	Multidimensional	Relational
Latency	low	high	Medium

CATEGORY	OLAP - online Analytical processing
Definition	db query mgmt syst
Data source	historical data
Method	Data Warehouse
Application	Data Mining / subject oriented
Normalized	OLAP Not normalized
Usage of data	Plan / problem solving
Task	multidimension
Purpose	analysis & decision making
Volume of data	large amount of data
Queries	slow
update	Not updated
process time	Complex
Types of users	Managed by CEO
operations	only read / rarely write
updates	Batch operation
Nature	focus customer
DB design	focus subject
Productivity	Improve efficiency

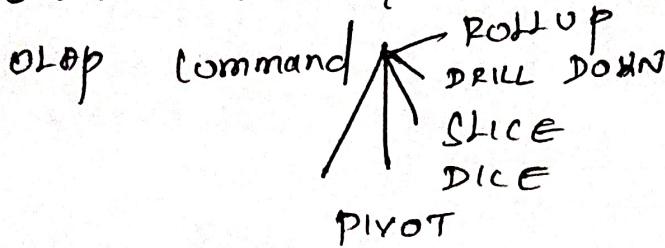
OLTP - online Transaction processing
db modifying systems
operational current data
DBMS
Business task / application oriented
3NF normalized
day to day operations
Business
Insert / update / delete
historical data
fast
maintained rigorously
fast
Managed by managers
Both read & write
Brief & quick
focus market
focus application
enhance productivity of users

TOPIC NAME:-

20"

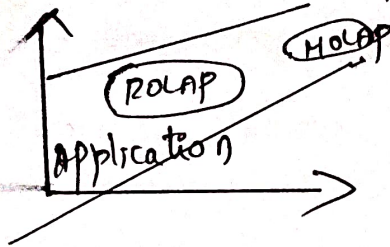
OLAP OPERATIONS - ONLINE ANALYTICAL PROCESSING

Execute user queries
SQL based methodology
describe location/time



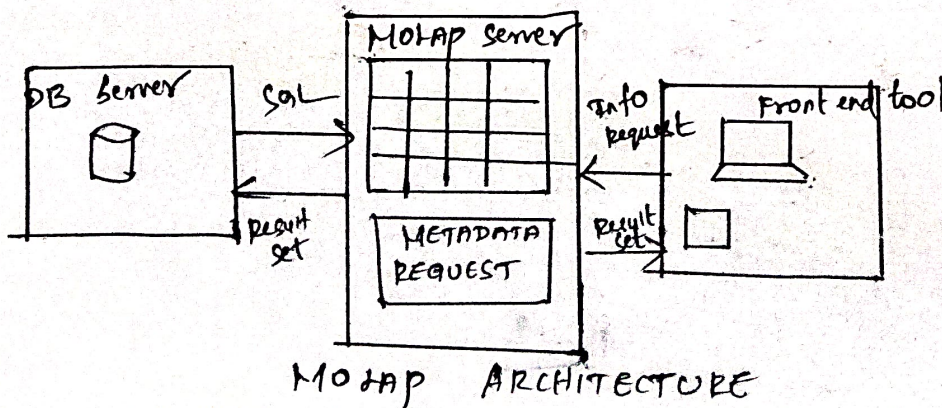
TYPES OF OLAP

OLAP tool multidimensional db
allow sophisticated users / analyze data
Multidimensional / complex views

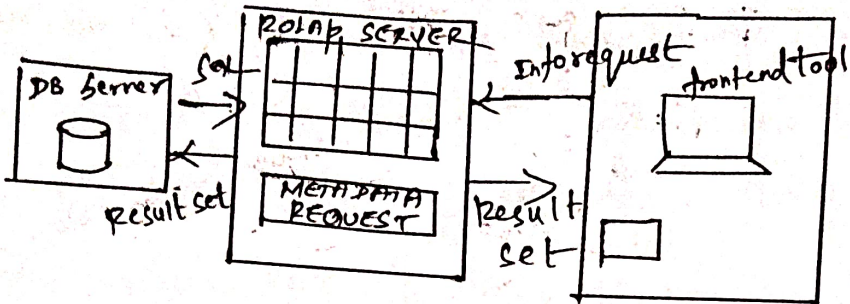


A) MOLAP :-

analyze data / aggregated form
predict complex analysis queries
Minimize disk space / perform time series analysis
Iterative / comprehensive time series
Access detail data / maintain RDBMS



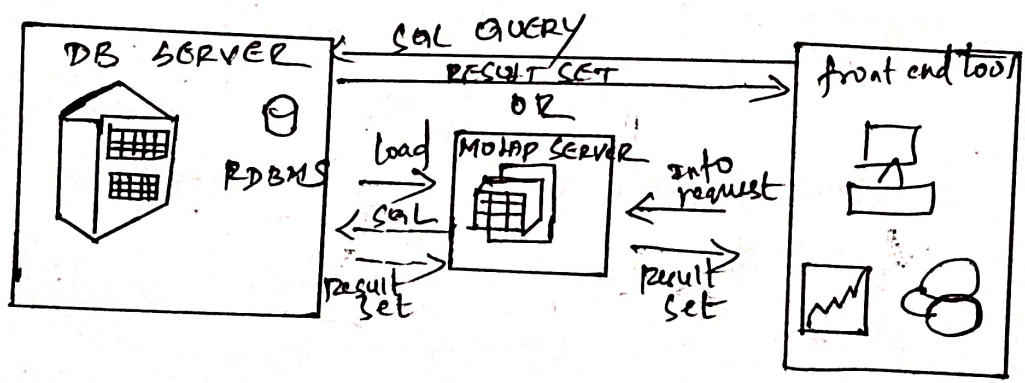
B) **ROLAP**:- fastest growing style / support RDBMS product
 create static multidimensional data structure



ROLAP ARCHITECTURE

C) **MQE** [MANAGED QUERY ENVIRONMENT]

provide "datacube" & "slice" & "dice" analysis
 store & maintained locally
 support complexity / handle user request
 Eg: sybase / allow flexibility



Hybrid / MQE - Architecture

— X —